

1

Introduction to Ridge Regression

This chapter reviews the developments of ridge regression, starting with the definition of ridge regression together with the covariance matrix. We discuss the multicollinearity problem and ridge notion and present the preliminary test and Stein-type estimators. In addition, we discuss the high-dimensional problem. In conclusion, we include detailed notes, references, and organization of the book.

1.1 Introduction

Consider the common multiple linear regression model with the vector of coefficients, $\beta = (\beta_1, \dots, \beta_p)^\top$ given by

$$Y = X\beta + \epsilon, \quad (1.1)$$

where $Y = (y_1, \dots, y_n)^\top$ is a vector of n responses, $X = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top$ is an $n \times p$ design matrix of rank p ($\leq n$), $\mathbf{x}_i \in \mathbb{R}^p$ is the vector of covariates, and ϵ is an n -vector of independently and identically distributed (i.i.d.) random variables (r.v.).

The least squares estimator (LSE) of β , denoted by $\tilde{\beta}_n$, can be obtained by minimizing the residual sum of squares (RSS), the convex optimization problem,

$$\min_{\beta} \{(Y - X\beta)^\top(Y - X\beta)\} = \min_{\beta} \{S(\beta)\},$$

where $S(\beta) = Y^\top Y - 2\beta^\top X^\top Y + \beta^\top X^\top X\beta$ is the RSS. Solving

$$\frac{\partial S(\beta)}{\partial \beta} = -2X^\top Y + 2X^\top X\beta = \mathbf{0}$$

with respect to (w.r.t.) β gives

$$\tilde{\beta}_n = (X^T X)^{-1} X^T Y. \quad (1.2)$$

Suppose that $\mathbb{E}(\epsilon) = 0$ and $\mathbb{E}(\epsilon\epsilon^T) = \sigma^2 \mathbf{I}_n$ for some $\sigma^2 \in \mathbb{R}^+$. Then, the variance–covariance matrix of LSE is given by

$$\text{Var}(\tilde{\beta}_n) = \sigma^2 (X^T X)^{-1}. \quad (1.3)$$

Now, we consider the canonical form of the multiple linear regression model to illustrate how large eigenvalues of the design matrix $X^T X$ may affect the efficiency of estimation.

Write the spectral decomposition of the positive definite design matrix $X^T X$ to get $X^T X = \Gamma \Lambda \Gamma^T$, where $\Gamma (p \times p)$ is a column orthogonal matrix of eigenvectors and $\Lambda = \text{Diag}(\lambda_1, \dots, \lambda_p)$, where $\lambda_j > 0, j = 1, \dots, p$ is the ordered eigenvalue matrix corresponding to $X^T X$. Then,

$$Y = T\xi + \epsilon, \quad T = X\Gamma, \quad \xi = \Gamma^T \beta. \quad (1.4)$$

The LSE of ξ has the form,

$$\begin{aligned} \tilde{\xi}_n &= (T^T T)^{-1} T^T Y \\ &= \Lambda^{-1} T^T Y. \end{aligned} \quad (1.5)$$

The variance–covariance matrix of $\tilde{\xi}_n$ is given by

$$\text{Var}(\tilde{\xi}_n) = \sigma^2 \Lambda^{-1} = \sigma^2 \begin{bmatrix} \frac{1}{\lambda_1} & 0 & \cdots & 0 \\ 0 & \frac{1}{\lambda_2} & 0 & 0 \\ \vdots & \cdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{1}{\lambda_p} \end{bmatrix}. \quad (1.6)$$

Summation of the diagonal elements of the variance–covariance matrix of $\tilde{\xi}_n$ is equal to $\text{tr}(\text{Var}(\tilde{\xi}_n)) = \sigma^2 \sum_{j=1}^p \lambda_j^{-1}$. Apparently, small eigenvalues inflate the total variance of estimate or energy of $X^T X$. Specifically, since the eigenvalues are ordered, if the first eigenvalue is small, it causes the variance to explode. If this happens, what must one do? In the following section, we consider this problem. Therefore, it is of interest to realize when the eigenvalues become small.

Before discussing this problem, a very primitive understanding is that if we enlarge the eigenvalues from λ_j to $\lambda_j + k$, for some positive value, say, k , then we can prevent the total variance from exploding. Of course, the amount of recovery depends on the correct choice of the parameter, k .

An artificial remedy is to have $\text{tr}(\text{Var}(\tilde{\xi}_n)) = \sigma^2 \sum_{j=1}^p (\lambda_j + k)^{-1}$ based on the variance matrix given by

$$\text{Var}(\tilde{\xi}_n) = \sigma^2 \begin{pmatrix} \frac{1}{\lambda_1 + k} & 0 & \cdots & 0 \\ 0 & \frac{1}{\lambda_2 + k} & 0 & 0 \\ \vdots & \cdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{1}{\lambda_p + k} \end{pmatrix} = \sigma^2 (\mathbf{\Lambda} + k\mathbf{I}_p)^{-1}, \quad k \in \mathbb{R}^+. \quad (1.7)$$

Replacing the eigenvector matrix $\mathbf{\Gamma}$ in (1.5) by this matrix (1.7), we get the $\tilde{\xi}_n = (\mathbf{\Lambda} + k\mathbf{I}_p)^{-1} \mathbf{T}^\top \mathbf{Y}$ and the variance as in (1.8).

$$\begin{aligned} \text{Var}(\tilde{\xi}_n) &= (\mathbf{\Lambda} + k\mathbf{I}_p)^{-1} \mathbf{T}^\top \text{Var}(\mathbf{Y}) \mathbf{T} (\mathbf{\Lambda} + k\mathbf{I}_p)^{-1} \\ &= \sigma^2 (\mathbf{\Lambda} + k\mathbf{I}_p)^{-1} \mathbf{\Lambda} (\mathbf{\Lambda} + k\mathbf{I}_p)^{-1}, \end{aligned} \quad (1.8)$$

which shows

$$\begin{aligned} \text{tr}(\text{Var}(\tilde{\xi}_n)) &= \sigma^2 \text{tr}((\mathbf{\Lambda} + k\mathbf{I}_p)^{-1} \mathbf{\Lambda} (\mathbf{\Lambda} + k\mathbf{I}_p)^{-1}) \\ &= \sigma^2 \sum_{j=1}^p \frac{\lambda_j}{(\lambda_j + k)^2} \end{aligned} \quad (1.9)$$

Further, we show that achieving the total variance of $\sigma^2 \sum_{j=1}^p \lambda_j / (\lambda_j + k)^2$ is the target.

1.1.1 Multicollinearity Problem

Multicollinearity or collinearity is the existence of near-linear relationships among the regressors, predictors, or input/exogenous variables. There are terms such as exact, complete and severe, or supercollinearity and moderate collinearity. Supercollinearity indicates that two (or multiple) covariates are linearly dependent, and moderate occurs when covariates are moderately correlated. In the complete collinearity case, the design matrix is not invertible. This case mostly occurs in a high-dimensional situation (e.g. microarray measure) in which the number of covariates (p) exceeds the number of samples (n).

Moderation occurs when the relationship between two variables depends on a third variable, namely, the moderator. This case mostly happens in structural equation modeling. Although moderate multicollinearity does not cause the mathematical problems of complete multicollinearity, it does affect the

interpretation of model parameter estimates. According to Montgomery et al. (2012), if there is no linear relationship between the regressors, they are said to be orthogonal.

Multicollinearity or ill-conditioning can create inaccurate estimates of the regression coefficients, inflate the standard errors of the regression coefficients, deflate the partial t -tests for the regression coefficients, give false and nonsignificant p -values, and degrade the predictability of the model. It also causes changes in the direction of signs of the coefficient estimates. According to Montgomery et al. (2012), there are five sources for multicollinearity: (i) data collection, (ii) physical constraints, (iii) overdefined model, (iv) model choice or specification, and (v) outliers.

There are many studies that well explain the problem of multicollinearity. Since theoretical aspects of ridge regression and related issues are our goal, we refer the reader to Montgomery et al. (2012) for illustrative examples and comprehensive study on the multicollinearity and diagnostic measures such as correlation matrix, eigen system analysis of $X^T X$, known as condition number, or variance decomposition proportion and variance inflation factor (VIF). To end this section, we consider a frequently used example in a ridge regression, namely, the Portland cement data introduced by Woods et al. (1932) from Najarian et al. (2013). This data set has been analyzed by many authors, e.g. Kaciranlar et al. (1999), Kibria (2003), and Arashi et al. (2015). We assemble the data as follows:

$$X = \begin{bmatrix} 7 & 26 & 6 & 60 \\ 1 & 29 & 15 & 52 \\ 11 & 56 & 8 & 20 \\ 11 & 31 & 8 & 47 \\ 7 & 52 & 6 & 33 \\ 11 & 55 & 9 & 22 \\ 3 & 71 & 17 & 6 \\ 1 & 31 & 22 & 44 \\ 2 & 54 & 18 & 22 \\ 21 & 47 & 4 & 26 \\ 1 & 40 & 23 & 34 \\ 11 & 66 & 9 & 12 \\ 10 & 68 & 8 & 12 \end{bmatrix}, \quad Y = \begin{bmatrix} 78.5 \\ 74.3 \\ 104.3 \\ 87.6 \\ 95.9 \\ 109.2 \\ 102.7 \\ 72.5 \\ 93.1 \\ 115.9 \\ 83.8 \\ 113.3 \\ 109.4 \end{bmatrix}. \quad (1.10)$$

The following (see Tables 1.1 and 1.2) is the partial output of linear regression fit to this data using the software SPSS, where we selected Enter as the method.

We display the VIF values to diagnose the multicollinearity problem in this data set. Values greater than 10 are a sign of multicollinearity. Many softwares can be used.

Table 1.1 Model fit indices for Portland cement data.

Model	R	R^2	Adjusted R^2	Standard error of the estimate
1	0.991	0.982	0.974	2.446

Table 1.2 Coefficient estimates for Portland cement data.

Coefficient	Standardized coefficient	t	Significance	VIF
(Constant)		0.891	0.399	
x_1	0.607	2.083	0.071	38.496
x_2	0.528	0.705	0.501	254.423
x_3	0.043	0.135	0.896	46.868
x_4	-0.160	-0.203	0.844	282.513

1.2 Ridge Regression Estimator: Ridge Notion

If the regression coefficients β_j s are unconstrained, then they can explode (become large); this results in high variance. Hence, in order to control the variance, one may regularize the regression coefficients and determine how large the coefficient grows. In other words, one may impose a constraint on them so as not to get unboundedly large or penalized large regression coefficients. One type of constraint is the ridge constraint given by $\sum_{j=1}^p \beta_j^2 \leq t$ for some positive value t . Hence, the minimization of the penalized residual sum of squares (PRSS) is equivalent to solving the following convex optimization problem,

$$\min_{\beta} \{(Y - X\beta)^\top (Y - X\beta)\} \quad \text{such that} \quad \sum_{j=1}^p \beta_j^2 \leq t, \quad (1.11)$$

for some positive value t .

In general, the PRSS is defined by

$$(Y - X\beta)^\top (Y - X\beta) + k\|\beta\|^2, \quad \|\beta\|^2 = \sum_{j=1}^p \beta_j^2. \quad (1.12)$$

Since the PRSS is a convex function w.r.t. β , it has a unique solution. Because of the ridge constraint, the solution is termed as the ridge regression estimator (RRE).

To derive the RRE, we solve the following convex optimization problem

$$\min_{\beta} \{(Y - X\beta)^\top (Y - X\beta) + k\|\beta\|^2\} = \min_{\beta} \{\text{PS}(\beta)\}, \quad (1.13)$$

where $\text{PS}(\beta) = Y^\top Y - 2\beta^\top X^\top Y + \beta^\top X^\top X\beta + k\beta^\top \beta$ is the PRSS. Solving

$$\frac{\partial \text{PS}(\beta)}{\partial \beta} = -2X^\top Y + 2X^\top X\beta + 2k\beta = \mathbf{0}$$

w.r.t. β gives the RRE,

$$\hat{\beta}_n^{\text{RR}}(k) = (X^\top X + kI_p)^{-1} X^\top Y. \quad (1.14)$$

Here, k is the shrinkage (tuning) parameter. Indeed, k tunes (controls) the size of the coefficients, and hence regularizes them. As $k \rightarrow 0$, the RRE simplifies to the LSE. Also, as $k \rightarrow \infty$, the RREs approach zero. Hence, the optimal shrinkage parameter k is of interest.

One must note that solving the optimization problem (1.13) is not the only way of yielding the RRE. It can also be obtained by solving a RSS of another data, say augmented data. To be specific, consider the following augmentation approach. Let

$$X^* = \begin{bmatrix} X \\ \sqrt{k}I_p \end{bmatrix}, \quad Y^* = \begin{bmatrix} Y \\ \mathbf{0} \end{bmatrix}. \quad (1.15)$$

Assume the following multiple linear model,

$$Y^* = X^* \beta + \epsilon^*, \quad (1.16)$$

where ϵ^* is an $(n+p)$ -vector of i.i.d. random variables. Then, the LSE of β is obtained as

$$\begin{aligned} \beta_n^* &= \min_{\beta} \{(Y^* - X^* \beta)^\top (Y^* - X^* \beta)\} \\ &= (X^{*\top} X^*)^{-1} X^{*\top} Y^* \\ &= (X^\top X + kI_p)^{-1} X^\top Y \\ &= \hat{\beta}_n^{\text{RR}}(k). \end{aligned} \quad (1.17)$$

Thus, the LSE of the augmented data is indeed the RRE of the normal data.

1.3 LSE vs. RRE

Indeed, the RRE is proportional to the LSE. Under the orthonormal case, i.e. $X^\top X = I_p$, the RRE simplifies to

$$\begin{aligned} \hat{\beta}_n^{\text{RR}}(k) &= (X^\top X + kI_p)^{-1} X^\top Y \\ &= (1+k)^{-1} I_p X^\top Y \\ &= (1+k)^{-1} (X^\top X)^{-1} X^\top Y \\ &= \frac{1}{1+k} \tilde{\beta}_n. \end{aligned} \quad (1.18)$$

Theorem 1.1 *The ridge penalty shrinks the eigenvalues of the design matrix.*

Proof: Write the singular value decomposition (SVD) of the design matrix X to get $X = \mathbf{U}\mathbf{\Lambda}\mathbf{V}^\top$, where \mathbf{U} ($n \times p$) and \mathbf{V} ($p \times p$) are column orthogonal matrices, and $\mathbf{\Lambda} = \text{Diag}(\lambda_1, \dots, \lambda_p)$ $\lambda_j > 0, j = 1, \dots, p$ is the eigenvalue matrix corresponding to $X^\top X$. Then,

$$\begin{aligned}\tilde{\beta}_n &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \\ &= (\mathbf{V}\mathbf{\Lambda}\mathbf{U}^\top \mathbf{U}\mathbf{\Lambda}\mathbf{V}^\top)^{-1} \mathbf{V}\mathbf{\Lambda}\mathbf{U}^\top \mathbf{Y} \\ &= (\mathbf{V}\mathbf{\Lambda}^2 \mathbf{V}^\top)^{-1} \mathbf{V}\mathbf{\Lambda}\mathbf{U}^\top \mathbf{Y} \\ &= \mathbf{V}\mathbf{\Lambda}^{-2} \mathbf{V}^\top \mathbf{V}\mathbf{\Lambda}\mathbf{U}^\top \mathbf{Y} \\ &= \mathbf{V}\{(\mathbf{\Lambda}^2)^{-1} \mathbf{\Lambda}\} \mathbf{U}^\top \mathbf{Y}.\end{aligned}\tag{1.19}$$

In a similar manner, one obtains

$$\begin{aligned}\hat{\beta}_n^{\text{RR}}(k) &= (\mathbf{X}^\top \mathbf{X} + k\mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{Y} \\ &= (\mathbf{V}\mathbf{\Lambda}\mathbf{U}^\top \mathbf{U}\mathbf{\Lambda}\mathbf{V}^\top + k\mathbf{I}_p)^{-1} \mathbf{V}\mathbf{\Lambda}\mathbf{U}^\top \mathbf{Y} \\ &= (\mathbf{V}\mathbf{\Lambda}^2 \mathbf{V}^\top + k\mathbf{V}\mathbf{V}^\top)^{-1} \mathbf{V}\mathbf{\Lambda}\mathbf{U}^\top \mathbf{Y} \\ &= \mathbf{V}(\mathbf{\Lambda}^2 + k\mathbf{I}_p)^{-1} \mathbf{V}^\top \mathbf{V}\mathbf{\Lambda}\mathbf{U}^\top \mathbf{Y} \\ &= \mathbf{V}\{(\mathbf{\Lambda}^2 + k\mathbf{I}_p)^{-1} \mathbf{\Lambda}\} \mathbf{U}^\top \mathbf{Y}.\end{aligned}\tag{1.20}$$

The difference between the two estimators w.r.t. the SVD is in the bracket terms. Let

$$\begin{aligned}\mathbf{\Lambda}^{-2} \mathbf{\Lambda} &= \mathbf{\Lambda}^{-1} = \text{Diag}\left(\frac{1}{\lambda_1}, \dots, \frac{1}{\lambda_p}\right), \\ (\mathbf{\Lambda}^2 + k\mathbf{I}_p)^{-1} \mathbf{\Lambda} &= \text{Diag}\left(\frac{\lambda_1}{\lambda_1^2 + k}, \dots, \frac{\lambda_p}{\lambda_p^2 + k}\right).\end{aligned}$$

Since $\lambda_j/(\lambda_j^2 + k) \leq \lambda_j^{-1}$, the ridge penalty shrinks the λ_j s. \square

1.4 Estimation of Ridge Parameter

We can observe from Eq. (1.18) that the RRE heavily depends on the ridge parameter k . Many authors at different times worked in this area of research and developed and proposed different estimators for k . They considered various models such as linear regression, Poisson regression, and logistic regression models. To mention a few, Hoerl and Kennard (1970), Hoerl et al. (1975), McDonald and Galarneau (1975), Lawless and Wang (1976), Dempster et al. (1977), Gibbons (1981), Kibria (2003), Khalaf and Shukur (2005), Alkhamisi and Shukur (2008), Muniz and Kibria (2009), Gruber et al. (2010), Muniz et al. (2012), Mansson et al. (2010), Hefnawy and Farag (2013), Aslam (2014), and Arashi and Valizadeh (2015), and Kibria and Banik (2016), among others.

1.5 Preliminary Test and Stein-Type Ridge Estimators

In previous sections, we discussed the notion of RRE and how it shrinks the elements of the ordinary LSE. Sometimes, it is needed to shrink the LSE to a subspace defined by $\mathbf{H}\boldsymbol{\beta} = \mathbf{h}$, where \mathbf{H} is a $q \times p$ known matrix of full row rank q ($q \leq p$) and \mathbf{h} is a q vector of known constants. It is also termed as constraint or restriction. Such a configuration of the subspace is frequently used in the design of experiments, known as contrasts. Therefore, sometimes shrinking is for two purposes. We refer to this as *double shrinking*.

In general, unlike the Bayesian paradigm, correctness of the prior information $\mathbf{H}\boldsymbol{\beta} = \mathbf{h}$ can be tested on the basis of samples through testing $\mathcal{H}_0 : \mathbf{H}\boldsymbol{\beta} = \mathbf{h}$ vs. a set of alternatives. Following Fisher's recipe, we use the non-sample information $\mathbf{H}\boldsymbol{\beta} = \mathbf{h}$; if based on the given sample, we accept \mathcal{H}_0 . In situations where this prior information is correct, an efficient estimator is the one which satisfies this restriction, called the *restricted estimator*.

To derive the restricted estimator under a multicollinear situation, satisfying the condition $\mathbf{H}\boldsymbol{\beta} = \mathbf{h}$, one solves the following convex optimization problem,

$$\min_{\boldsymbol{\beta}} \{PS_{\lambda}(\boldsymbol{\beta})\}, \quad PS_{\lambda}(\boldsymbol{\beta}) = (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^{\top}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + k\|\boldsymbol{\beta}\|^2 + \boldsymbol{\lambda}^{\top}(\mathbf{H}\boldsymbol{\beta} - \mathbf{h}), \quad (1.21)$$

where $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_q)^{\top}$ is the vector of Lagrangian multipliers. Grob (2003) proposed the restricted RRE, under a multicollinear situation, by correcting the restricted RRE of Sarkar (1992).

In our case, we consider prior information with the form $\boldsymbol{\beta} = \mathbf{0}$, which is a test used for checking goodness of fit. Here, the restricted RRE is simply given by $\hat{\boldsymbol{\beta}}_n(k) = \mathbf{0}$, where $\mathbf{0}$ is the restricted estimator of $\boldsymbol{\beta}$. Therefore, one uses $\hat{\boldsymbol{\beta}}_n^{\text{RR}}(k)$ if $\mathcal{H}_0 : \boldsymbol{\beta} = \mathbf{0}$ rejects and $\hat{\boldsymbol{\beta}}_n^{\text{RR(R)}}(k)$ if $\mathcal{H}_0 : \boldsymbol{\beta} = \mathbf{0}$ accepts. Combining the information existing in both estimators, one may follow an approach by Bancroft (1964), to propose the preliminary test RRE given by

$$\begin{aligned} \hat{\boldsymbol{\beta}}_n^{\text{RR(PT)}}(k, \alpha) &= \begin{cases} \hat{\boldsymbol{\beta}}_n^{\text{RR}}(k); & \mathcal{H}_0 \text{ is rejected} \\ \hat{\boldsymbol{\beta}}_n^{\text{RR(R)}}(k); & \mathcal{H}_0 \text{ is accepted} \end{cases} \\ &= \hat{\boldsymbol{\beta}}_n^{\text{RR}}(k)I(\mathcal{H}_0 \text{ is rejected}) + \hat{\boldsymbol{\beta}}_n^{\text{RR(R)}}(k)I(\mathcal{H}_0 \text{ is accepted}) \\ &= \hat{\boldsymbol{\beta}}_n^{\text{RR}}(k)I(\mathcal{L}_n > F_{p,m}(\alpha)) + \hat{\boldsymbol{\beta}}_n^{\text{RR(R)}}(k)I(\mathcal{L}_n \leq F_{p,m}(\alpha)) \\ &= \hat{\boldsymbol{\beta}}_n^{\text{RR}}(k) - \hat{\boldsymbol{\beta}}_n^{\text{RR}}(k)I(\mathcal{L}_n \leq F_{p,m}(\alpha)), \end{aligned} \quad (1.22)$$

where $I(A)$ is the indicator function of the set A and \mathcal{L}_n is the test statistic for testing $\mathcal{H}_0 : \boldsymbol{\beta} = \mathbf{0}$, and $F_{q,m}(\alpha)$ is the upper α -level critical value from the F -distribution with (q, m) degrees of freedom (D.F.) See Judge and Bock (1978) and Saleh (2006) for the test statistic and details.

After some algebra, it can be shown that

$$\hat{\beta}_n^{\text{RR(R)}}(k) = \mathbf{R}(k)\hat{\beta}_n, \quad \hat{\beta}_n = \mathbf{0}, \quad (1.23)$$

where

$$\mathbf{R}(k) = (\mathbf{I}_p + k(\mathbf{X}^\top \mathbf{X})^{-1})^{-1}.$$

Then, it is easy to show that

$$\hat{\beta}_n^{\text{RR(PT)}}(k, \alpha) = \mathbf{R}(k)\hat{\beta}_n^{\text{PT}}, \quad \hat{\beta}_n^{\text{PT}} = \tilde{\beta}_n - \tilde{\beta}_n I(\mathcal{L}_n \leq F_{p,m}(\alpha)). \quad (1.24)$$

The preliminary test RRE is discrete in nature and heavily dependent on α , the level of significance. Hence, a continuous α -free estimator is desired. Following James and Stein (1961), the Stein-type RRE is given by

$$\hat{\beta}_n^{\text{RR(S)}}(k) = \mathbf{R}(k)\hat{\beta}_n^{\text{S}}; \quad \hat{\beta}_n^{\text{S}} = \tilde{\beta}_n - d\mathcal{L}_n^{-1}\tilde{\beta}_n, \quad (1.25)$$

for some $d > 0$. It is shown in Saleh (2006) that the optimal choice of d is equal to $m(p+2)/p(m+2)$, $m = n - p$.

Note that if $\mathcal{L}_n \rightarrow \infty$, then $\hat{\beta}_n^{\text{RR(S)}}(k) = \hat{\beta}_n^{\text{RR}}(k)$, which matches with the estimator $\hat{\beta}_n^{\text{RR(PT)}}(k)$. However, if $\mathcal{L}_n \rightarrow 0$, then $\hat{\beta}_n^{\text{RR(S)}}(k)$ becomes a negative estimator. In order to obtain the true value, one has to restrict ($\mathcal{L}_n > d$). Hence, one may define the positive-rule Stein-type RRE given by

$$\hat{\beta}_n^{\text{RR(S+)}}(k) = \mathbf{R}(k)\hat{\beta}_n^{\text{S+}}; \quad \hat{\beta}_n^{\text{S+}} = \hat{\beta}_n^{\text{S}} - (1 - d\mathcal{L}_n^{-1})I(\mathcal{L}_n \leq d)\tilde{\beta}_n. \quad (1.26)$$

Although these shrinkage RREs are biased, they outperform the RRE, $\hat{\beta}_n^{\text{RR}}(k)$ in the mean squared error (MSE) sense. One must note that this superiority is not uniform over the parameter space \mathbb{R}^p . According to how much β deviates from the origin, the superiority changes. We refer to Saleh (2006) for an extensive overview on this topic and statistical properties of the shrinkage estimators. This type of estimator will be studied in more detail and compared with the RRE in the following chapters.

1.6 High-Dimensional Setting

In high-dimensional analysis, the number of variables p is greater than the number of observations, n . In such situations, $\mathbf{X}^\top \mathbf{X}$ is not invertible and, hence, to derive the LSE of β , one may use the generalized inverse of $\mathbf{X}^\top \mathbf{X}$, which does not give a unique solution. However, the RRE can be obtained, regardless of the relationship between p and n . In this section, we briefly mention some related endeavors in manipulating the RRE to adapt with high-dimensional setting. However, in Chapter 11, we discuss the growing dimension, i.e. $p \rightarrow \infty$ when the sample size is fixed.

Wang et al. (2016) used this fact and proposed a new approach in high-dimensional regression by considering

$$(X^T X)^{-1} X^T Y = \lim_{k \rightarrow 0} (X^T X + k I_p)^{-1} X^T Y. \quad (1.27)$$

They used an orthogonal projection of β onto the row space of X and proposed the following high-dimensional version of the LSE for dimension reduction:

$$\hat{\beta}_n^{\text{HD}} = \lim_{k \rightarrow 0} X^T (X X^T + k I_n)^{-1} Y = X^T (X X^T)^{-1} Y. \quad (1.28)$$

They applied the following identity to obtain the estimator:

$$(X^T X + k I_p)^{-1} X^T Y = X^T (X X^T + k I_n)^{-1} Y \quad (1.29)$$

for every $n, p, k > 0$. Buhlmann et al. (2014) also used the projection of β onto the row space of X and developed a bias correction in the RRE to propose a bias-corrected RRE for the high-dimensional setting. Shao and Deng (2012) considered the RRE of the projection vector and proposed to threshold the RRE when the projection vector is sparse, in the sense that many of its components are small.

Specifically, write the SVD of X as $X = U \Lambda V^T$ where $U(n \times p)$ and $V(p \times p)$ are column orthogonal matrices, and let $\Lambda = \text{Diag}(\lambda_1, \dots, \lambda_p)$, $\lambda_j > 0$, $j = 1, \dots, p$ be the eigenvalue matrix corresponding to $X^T X$. Let

$$\theta = X^T (X X^T)^{-1} X \beta = V V^T \beta, \quad (X X^T)^{-1} = U \Lambda^{-2} U^T, \quad (1.30)$$

where θ is the projection of β onto the row space of X .

Since $X X^T (X X^T)^{-1} X = X$, we have $X \theta = X \beta$. Obviously, $X \theta = X \beta$ yields $Y = X \theta + \epsilon$ as the underlying model instead of (1.1). Then, the RRE of the projection is simply

$$\hat{\theta}_n^{\text{RR}}(k_n) = (X^T X + k_n I_p)^{-1} X^T Y, \quad (1.31)$$

where $k_n > 0$ is an appropriately chosen regularization parameter.

Shao and Deng (2012) proved that the RRE $\hat{\theta}_n(k_n)$ is consistent for the estimation of any linear combination of θ ; however, it is not L_2 consistent, i.e. $n^{-1} \mathbb{E} \|X \hat{\theta}_n(k_n) - X \theta\|^2$ may not converge to zero. Then, they proposed an L_2 -consistent RRE by thresholding.

Another approach to deal with high-dimensional data in the case $p > n$ is to partition the vector of regression parameters to main and nuisance effects, where we have p_1 (say) less than n , main and active coefficients that must be estimated and $p_2 = p - p_1$, redundant coefficients that can be eliminated from the analysis by a variable selection method. In fact, we assume some level of sparsity in our inference. Hence, we partition $\beta = (\beta_1^T, \beta_2^T)^T$ in such a way in which β_1 is a p_1 vector of main effects and β_2 is a p_2 vector of nuisance effects. We further assume that $p_1 < n$. Accordingly, we have the partition $X = (X_1, X_2)$ for the design matrix. Then, the multiple linear model (1.1) is rewritten as

$$Y = X_1 \beta_1 + X_2 \beta_2 + \epsilon. \quad (1.32)$$

This model may be termed as the full model. As soon as a variable selection method is used and p_1 main effects are selected, the sub model has the form

$$Y = X_1\beta_1 + \epsilon. \quad (1.33)$$

Then, the interest is to estimate the main effects, β_1 (this technique is discussed in more detail in Chapter 11). Therefore, one needs a variable selection method. Unlike the RRE that only shrinks the coefficients, there are many estimators which simultaneously shrink and select variables in large p and small n problems. The most well-known one is the least absolute shrinkage and selection operator (LASSO) proposed by Tibshirani (1996). He suggested using an absolute-type constraint with form $\sum_{j=1}^p |\beta_j| \leq t$ for some positive values t in the minimization of the PRSS rather than the ridge constraint. Specifically, he defined the PRSS by

$$(Y - X\beta)^\top(Y - X\beta) + k\|\beta\|^\top \mathbf{1}_p, \quad \|\beta\|^\top \mathbf{1}_p = \sum_{j=1}^p |\beta_j|, \quad (1.34)$$

where $\|\beta\| = (|\beta_1|, \dots, |\beta_p|)^\top$ and $\mathbf{1}_p$ is an n -tuple of 1's. The LASSO estimator, denoted by $\hat{\beta}_n^{\text{LASSO}}$, is the solution to the following convex optimization problem,

$$\min_{\beta} (Y - X\beta)^\top(Y - X\beta) + k\|\beta\|^\top \mathbf{1}_p. \quad (1.35)$$

To see some recent related endeavors in the context of ridge regression, we refer to Yuzbasi and Ahmed (2015) and Aydin et al. (2016) among others.

In the following section, some of the most recent important references about ridge regression and related topics are listed for more studies.

1.7 Notes and References

The first paper on ridge analysis was by Hoerl (1962); however, the first paper on multicollinearity appeared five years later, roughly speaking, by Farrar and Glauber (1967). Marquardt and Snee (1975) reviewed the theory of ridge regression and its relation to generalized inverse regression. Their study includes several illustrative examples about ridge regression. For the geometry of multicollinearity, see Akdeniz and Ozturk (1981). We also suggest that Gunst (1983) and Sakallioglu and Akdeniz (1998) not be missed. Gruber (1998) in his monograph motivates the need for using ridge regression and allocated a large portion to the analysis of ridge regression and its generalizations. For historical survey up to 1998, we refer to Gruber (1998).

Beginning from 2000, a comprehensive study in ridge regression is the work of Ozturk and Akdeniz (2000), where the authors provide some solutions for ill-posed inverse problems. Wan (2002) incorporated measure of goodness of fit in evaluating the RRE and proposed a feasible generalized RRE. Kibria

(2003) gave a comprehensive analysis about the estimation of ridge parameter k for the linear regression model. For application of ridge regression in agriculture, see Jamal and Rind (2007). Maronna (2011) proposed an RRE based on repeated M-estimation in robust regression. Saleh et al. (2014) extensively studied the performance of preliminary test and Stein-type ridge estimators in the multivariate- t regression model. Huang et al. (2016) defined a weighted VIF for collinearity diagnostic in generalized linear models.

Arashi et al. (2017) studied the performance of several ridge parameter estimators in a restricted ridge regression model with stochastic constraints. Asar et al. (2017) defined a restricted RRE in the logistic regression model and derived its statistical properties. Roozbeh and Arashi (2016a) developed a new ridge estimator in partial linear models. Roozbeh and Arashi (2016b) used difference methodology to study the performance of an RRE in a partial linear model. Arashi and Valizadeh (2015) compared several estimators for estimating the biasing parameter in the study of partial linear models in the presence of multicollinearity. Roozbeh and Arashi (2013) proposed a feasible RRE in partial linear models and studied its properties in details. Roozbeh et al. (2012) developed RREs in seemingly partial linear models. Recently, Chandrasekhar et al. (2016) proposed the concept of partial ridge regression, which involves selectively adjusting the ridge constants associated with highly collinear variables to control instability in the variances of coefficient estimates. Norouzirad and Arashi (2017) developed shrinkage ridge estimators in the context of robust regression. Fallah et al. (2017) studied the asymptotic performance of a general form of shrinkage ridge estimator. Recently, Norouzirad et al. (2017) proposed improved robust ridge M-estimators and studied their asymptotic behavior.

1.8 Organization of the Book

This book has 12 chapters including this one. In this light we consider the chapter with location and simple linear model first. Chapter 1 presents an introduction to ridge regression and different aspects of it, stressing the multicollinearity problem and its application to high-dimensional problems. Chapter 2 considers simple linear and location models, and provides theoretical developments. Chapters 3 and 4 deal with the analysis of variance (ANOVA) model and seemingly unrelated simple linear models, respectively. Chapter 5 considers ridge regression and LASSO for multiple regression together with preliminary test and Stein-type estimator and a comparison thereof when the design matrix is non-orthogonal. Chapter 6 considers the RRE and its relation to LASSO. Further, we study in detail the properties of the preliminary test and Stein-type estimator with low dimension. In Chapter 7,

we cover the partially linear model and the properties of LASSO, ridge, preliminary test, and the Stein-type estimators. Chapter 8 contains the discussion on the logistic regression model and the related estimators of diverse kinds as described in earlier chapters. Chapter 9 discusses the multiple regression model with autoregressive errors. In Chapter 10, we provide a comparative study of LASSO, ridge, preliminary test, and Stein-type estimators using rank-based theory. In Chapter 11, we discuss the estimation of parameters of a regression model with high dimensions. Finally, we conclude the book with Chapter 12 to illustrate recent applications of ridge, LASSO, and logistic regression to neural networks and big data analysis.

Problems

- 1.1 Derive the estimates given in (1.12) and (1.13).
- 1.2 Find the optimum value $k > 0$ for which the MSE of RRE in (1.18) becomes the smallest.
- 1.3 Derive the bias and MSE functions for the shrinkage RREs given by (1.24)–(1.26).
- 1.4 Verify that

$$\hat{\beta}_n^{\text{RR}}(k) = (\mathbf{I}_p + k(\mathbf{X}^\top \mathbf{X})^{-1})^{-1} \tilde{\beta}_n.$$
- 1.5 Verify the identity (1.29).
- 1.6 Find the solution of the optimization problem (1.35) for the orthogonal case $\mathbf{X}^\top \mathbf{X} = n\mathbf{I}_p$.

